

1 A comprehensive LFQ benchmark dataset on modern day acquisition strategies 2 in proteomics.

3 Bart Van Puyvelde¹, Simon Daled¹, Sander Willems², Ralf Gabriels^{3,4}, Anne Gonzalez de Peredo⁵, Karima
4 Chaoui⁵, Emmanuelle Mouton-Barbosa⁵, David Bouyssié⁵, Kurt Boonen^{6,7}, Christopher J. Hughes⁸, Lee A.
5 Gethings⁸, Yasset Perez-Riverol⁹, Odile Schiltz⁵, Lennart Martens^{3,4}, Dieter Deforce¹, Maarten
6 Dhaenens^{1*}

7 1. ProGenTomics, Laboratory of Pharmaceutical Biotechnology, Ghent University, Ghent, Belgium
8 2. Department of Proteomics and Signal Transduction, Max Planck Institute of Biochemistry, Martinsried, Germany
9 3. VIB-UGent Center for Medical Biotechnology, VIB, 9000 Ghent, Belgium
10 4. Department of Biomolecular Medicine, Ghent University, 9000 Ghent, Belgium
11 5. Institut de Pharmacologie et de Biologie Structurale (IPBS), Université de Toulouse, CNRS, UPS, Toulouse, France
12 6. VITO Health, Mol, Belgium
13 7. Centre for Proteomics, University of Antwerpen, Antwerp, Belgium
14 8. Waters Corporation, Wilmslow, UK
15 9. European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, UK

16
17 *Corresponding author: Maarten Dhaenens (maarten.dhaenens@ugent.be)

18 **In the last decade, a revolution in liquid chromatography-mass spectrometry (LC-MS) based proteomics**
19 **was unfolded with the introduction of dozens of novel instruments that incorporate additional data**
20 **dimensions through innovative acquisition methodologies, in turn inspiring specialized data analysis**
21 **pipelines. Simultaneously, a growing number of proteomics datasets have been made publicly available**
22 **through data repositories such as ProteomeXchange, Zenodo and Skyline Panorama. However, developing**
23 **algorithms to mine this data and assessing the performance on different platforms is currently hampered**
24 **by the lack of a single benchmark experimental design. Therefore, we acquired a hybrid proteome mixture**
25 **on different instrument platforms and in all currently available families of data acquisition. Here, we**
26 **present a comprehensive Data-Dependent and Data-Independent Acquisition (DDA/DIA) dataset acquired**
27 **using several of the most commonly used current day instrumental platforms. The dataset consists of over**
28 **700 LC-MS runs, including adequate replicates allowing robust statistics and covering over nearly 10**
29 **different data formats, including scanning quadrupole and ion mobility enabled acquisitions. Datasets are**
30 **available via ProteomeXchange (PXD028735).**

31 **Background & Summary**

32 Hypothesis-driven biochemical assays have been the foundation of molecular biology for well over a
33 century, with great success. However, the lack of a more holistic view on the biomolecular complexity
34 requires trial-and-error experimentation. Therefore, the past few decades were characterized by a shift
35 towards an experimental design wherein a broader biomolecular perspective of the system is first
36 generated in order to contextualize the hypothesis and the targeted biochemical assays beforehand.
37 These “omics” approaches were enabled by two technical revolutions, i.e. the sequencing of nucleotides
38 and the accurate mass measurement of biomolecules by mass spectrometry (MS).

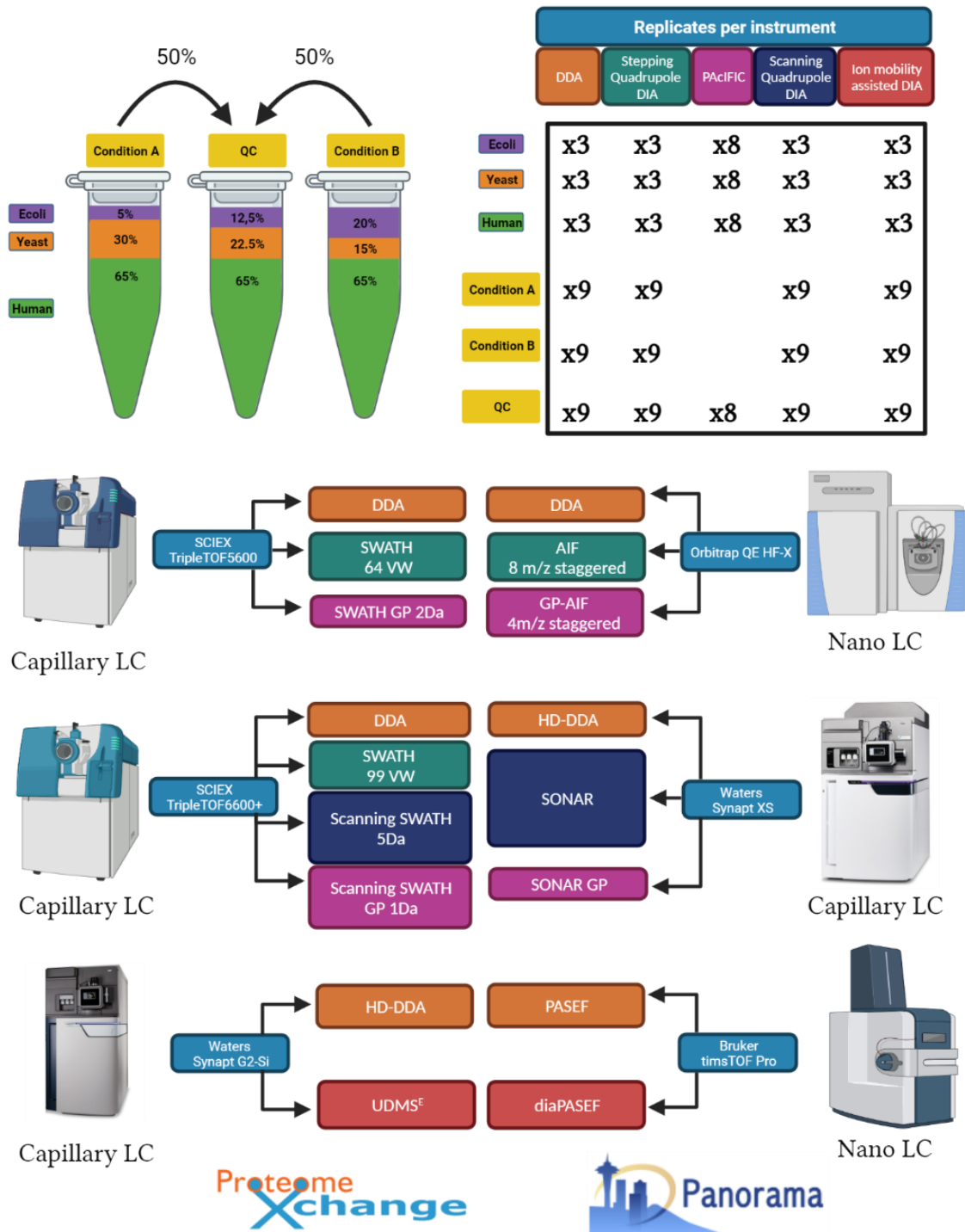
39 In its barest form, the output from an MS instrument is merely a list of m/z 's with intensities measured
40 at very precise moments in time. However, MS is quickly evolving towards capturing the full complexity
41 of a biological sample. To this end, not only the accuracy of instruments has improved greatly, they now
42 also incorporate analytical techniques that select or separate analytes based on other physico-chemical
43 properties. In proteomics nowadays, a mass spectrometer thus rarely only measures the m/z coordinate
44 and intensity of (fragment) ions. The ion coordinates are mostly supplemented with precursor m/z ,
45 retention time (t_R) and/or ion mobility coordinates, depending on acquisition strategy. This creates a
46 multidimensional data matrix that captures the complexity of the sample to an unprecedented depth ¹.

47 The field of mass spectrometry has diversified greatly, driven by a fast sequence of innovations from
48 many vendors. Instrumental engineering has allowed to manipulate ions in countless of different ways,
49 including different ionization methodologies, fragmentation techniques, multipoles, time-of-flight
50 tubes, ion mobility separation devices and trap designs, including the now very dominant orbitrap.
51 Especially the way in which these different ion manipulations are combined has ballooned the number
52 of different acquisition strategies available to the end user today. Irrespectively however, all these
53 instrumental and strategic innovations are futile if no data analysis pipeline is available to translate the
54 data back into biology. For bottom-up proteomics this implies reconstructing the peptide backbone
55 sequences from their fragment ions because the latter encompasses the specificity for identifying the
56 hundreds of millions of different protein sequences that make up the biotic world.

57 Conventionally, MS instruments have been operated using data dependent acquisition (DDA) wherein
58 the data from a precursor scan at low energy is used to pinpoint potentially interesting analytes which
59 are then sequentially selected for fragmentation at high energy. These fragmentation spectra can then
60 be identified by a plethora of different algorithms ^{2,3}. Data-independent acquisition (DIA) however, is
61 the more intuitive way of analyzing a sample, because it captures all (fragment) ions at an equal pace
62 without any instrumental bias. Yet, interpreting such complex data matrix has proven difficult and an
63 additional separation dimension, such as ion mobility, was added to increase the discriminating power
64 ⁴⁻⁸. Alternatively, configuring a quadrupole to sequentially scan the entire mass range - while still
65 operating “data independent” - alleviates the complexity of the resulting fragmentation spectra even
66 more ^{9,10}. This has opened up the way for the many different spectrum-centric and peptide-centric data
67 analysis strategies available today ¹¹⁻¹⁶. The latest reduction in complexity or “chimericity” of DIA spectra
68 encompasses continuously scanning the quadrupole as is done with SONAR ¹⁷ and Scanning SWATH ¹⁸
69 and combining quadrupole selection and ion mobility separation, as is done with diaPASEF ¹⁹.
70 Unsurprisingly, machine learning is taking center stage in mining the various resulting data architectures
71 ²⁰⁻²⁷.

72 Here, we created a comprehensive dataset on a single benchmark experimental design adapted from
73 Navarro et al. ²⁸. It contains a ground truth that serves as a quality control for bioinformatics algorithm
74 validation. This sample was acquired in adequate replicates on many of the current day instrumental
75 platforms – partially in nano flow LC and partially in capillary flow LC - by most of the available acquisition
76 strategy families, covering all commonly measured ion coordinates (**Figure 1**). Far from being complete,
77 it still is the most comprehensive repository of its kind for algorithmic development and validation, both
78 at the level of identification and quantification. Instead of being yet another way of attaining the highest
79 number of identified peptides, we hope it to become a resource for compatibility assessments and data
80 analysis quality control. Above all, it is a snapshot of current day completeness of our digital image of
81 the protein world.

82



83

84 **Figure.1 Schematic overview of the different acquisition strategies/instruments applied in the study.** A
 85 comprehensive LC-MS/MS dataset was generated using samples composed of commercial Human,
 86 Yeast and E.coli full proteome digests. Two hybrid proteome samples A and B containing known
 87 quantities of Human, Yeast and E.coli tryptic peptides, as described by Navarro et al. were prepared in
 88 three consecutive times to include handling variability. Additionally, a QC sample was created by mixing
 89 one sixth of each of the six master batches (65% w/w Human, 22.5% w/w Yeast and 12.5% w/w E.coli).
 90 These commercial lysates were measured individually and as triple hybrid proteome mixtures each in
 91 triplicate using DDA and DIA acquisition methodologies available on six LC-MS/MS platforms, i.e. SCIEX
 92 TripleTOF5600 and TripleTOF 6600+, Thermo Orbitrap QE HF-X, Waters Synapt G2-Si and Synapt XS and
 93 Bruker timsTOF Pro. The complete dataset was made publicly available to the proteomics community

94 through ProteomeXchange with dataset identifier: PXD028735. In addition, a system suitability
95 workflow (AutoQC) was incorporated on each instrument using commercial E.coli lysate digest which
96 were acquired at multiple timepoints throughout each sample batch. The AutoQC data was
97 automatically imported in Skyline and uploaded to the Panorama AutoQC server using AutoQC loader,
98 enabling system suitability assessment of each LC-MS/MS system used in the dataset.

99 **Methods**

100 **Sample preparation**

101 Mass spectrometry-compatible Human (P/N: V6951) and Yeast (P/N: V7461) protein digest extracts
102 were purchased from Promega (Madison, Wisconsin, United States). Lyophilised MassPrep E.coli digest
103 standard (P/N:186003196) was purchased from Waters Corporation (Milford, Massachusetts, United
104 States). The extracts were reduced with dithiothreitol (DTT), alkylated with iodoacetamide (IAA) and
105 digested with sequencing grade Trypsin(-Lys C) by the respective manufacturers. The digested protein
106 extracts were reconstituted in a mixture of 0.1% Formic acid (FA) in water (Biosolve B.V, Valkenswaard,
107 The Netherlands) and spiked with iRT peptides (Biognosys, Schlieren, Switzerland) at a ratio of 1:20 v/v.
108 Two master samples A and B were created similar to Navarro et al., each in triplicate, as shown in **Figure**
109 **1**. Sample A was prepared by mixing Human, Yeast and E.coli at 65% 30% and 20% weight for weight
110 (w/w), respectively. Sample B was prepared by mixing Human, Yeast and E.coli protein digests at 65%,
111 15%, 20% w/w, respectively. The resulting samples have logarithmic fold changes (log₂FCs) of 0,-1 and
112 2 for respectively Human, Yeast and E.coli. One sixth of each of the triplicate master batches of A and B
113 were mixed to create a QC sample, containing 65% w/w Human, 22.5% w/w Yeast and 12.5% w/w E.coli.

114 **LC-MS/MS**

115 In this section, a detailed description of the different LC-MS/MS parameters is given for each LC-MS/MS
116 instrumental setup applied to generate this comprehensive dataset. All instruments were operated
117 according to the lab's best practice, i.e. not necessarily the best attainable, but rather most realistic data
118 quality. Sample load was chosen based on LC setup (nano flow = 1 µg on column vs capillary flow = 5 µg
119 on column) and instrument sensitivity. Thus, differences in absolute number of identified peptides and
120 proteins can be attributed to sample load, LC flow rate, MS instrumentation, operator's choices and
121 search algorithmic compatibility; direct conclusions on MS instrument performance can therefore not
122 be drawn from this dataset.

123 As a rule of thumb, data-dependent acquisition (DDA) methods use high energy fragmentation spectra
124 (MS₂) of narrow mass selections for identification and use the area under the curve of the precursor
125 (MS₁) for quantification. Therefore, a cycle time needs to be attained wherein enough datapoints across
126 the precursor elution peak are sampled for accurate quantification. In most data-independent
127 acquisition (DIA) strategies, a broader precursor selection window is used and both identification and
128 quantification can be done at the fragment level, taken that the cycle time for both MS₁ and MS₂ is
129 adapted to the LC gradient. Finally, Precursor Acquisition Independent From Ion Count (PAcIFIC) is a
130 method that is acquired solely to extend the size of the peptide library for detecting peptides in DIA
131 data and is therefore not strictly dependent on the cycle time. Of note, by scanning the quadrupole
132 instead of acquiring different mass windows separately, acquisition strategies like SONAR and Scanning
133 SWATH create an additional dimension in the data matrix, akin to how ion mobility separation is
134 perceived. Since these are similar to PAcIFIC acquisition, we also acquired gas phase (GP) fractions in
135 SONAR and Scanning SWATH for library building, i.e. with no emphasis on cycle time.

136 **1) SCIEX TripleTOF 5600 (Capillary flow LC)**

137 A TripleTOF 5600 mass spectrometer (Sciex, Concord, Ontario, Canada) fitted with a Duospray ion
138 source operating in positive ion mode, was coupled to an Eksigent NanoLC 400 HPLC system (Eksigent,
139 Dublin, CA). 5 μ L of each sample was loaded at 5 μ L/min with 0.1% Trifluoroacetic acid (TFA) in water
140 and trapped on a YMC TriArt C18 guard column (id 500 μ m, length 5mm, particle size 3 μ m) for 5
141 minutes. Peptides were separated on a microLC YMC TriArt C18 column (id 300 μ m, length 15 cm,
142 particle size 3 μ m) maintained at 55°C at a flow rate of 5 μ L/min by means of trap-elute injection. Mobile
143 phase A consisted of UPLC-grade water with 0.1% (v/v) FA and 3% (v/v) DMSO, and mobile phase B
144 consisted of UPLC-grade ACN with 0.1% (v/v) FA. Peptide elution was performed at 5 μ L/min using the
145 following gradient: i) 2% to 30% mobile phase B in 120 min, ii) ramp to 90% mobile phase B in 1 min.
146 The washing step at 90% mobile phase B lasted 4 min and was followed by an equilibration step at 2%
147 mobile phase B (starting conditions) for 10 min. Ion source parameters were set to 5.5 kV for the ion
148 spray voltage, 30 psi for the curtain gas, 13 psi for the nebulizer gas and 80°C as source temperature.

149 a. Data-Dependent Acquisition

150 For DDA (a cycle time of 3.5 s), MS1 spectra were collected between 399-1200 m/z for 500 ms. The 20
151 most intense precursors ions with charge states 2-5 that exceeded 250 counts per second were selected
152 for fragmentation, and the corresponding fragmentation MS2 spectra were collected between 50-2000
153 m/z for 151 ms. After the fragmentation event, the precursor ions were dynamically excluded from
154 reselection for 20 s.

155 b. PAcIFIC

156 For PAcIFIC (a cycle time of 4 s), the TripleTOF5600 was configured to acquire eight gas phase
157 fractionated acquisitions with isolation windows of 4 m/z using an overlapping window pattern from
158 narrow mass ranges, as described by Searle et al (i.e., 396.43 – 502.48; 496.48 – 602.52; 596.52 –
159 702.57; 696.57 – 802.61; 796.61 – 902.66; 896.6 – 1002.70; 996.70 – 1102.75; 1096.75 – 1202.80)²⁹.
160 See Supplementary Data for the actual windowing scheme. MS2 spectra were collected in high-
161 sensitivity mode from 360-1460 m/z, for 75 ms. An MS1 survey scan was recorded per cycle from 360-
162 1460 m/z for 50ms.

163 c. SWATH 64 variable windows

164 For SWATH (a cycle time of 3.4 s), a 64 variable window acquisition scheme as described by Navarro et
165 al. was used for all samples²⁸. Briefly, SWATH MS2 spectra were collected in high-sensitivity mode from
166 50-2000 m/z, for 50 ms. Before each SWATH MS cycle an additional MS1 survey scan in high sensitivity
167 mode from 400-1200 m/z was recorded for 150 ms.

168 **2) SCIEX TripleTOF 6600+ (Capillary flow LC)**

169 A TripleTOF 6600+ mass spectrometer (Sciex, Concord, Ontario, Canada) fitted with an Optiflow ion
170 source operating in positive ion mode, was coupled to an Eksigent NanoLC 425 HPLC system (Eksigent,
171 Dublin, CA). 5 μ L of each sample was loaded at 5 μ L/min with 0.1% FA in water by means of direct
172 injection. Peptides were separated on a Phenomenex Luna Omega Polar C18 column (150 x 0.3 mm,
173 particle size 3 μ m) at a column temperature of 30°C. Mobile phase A consisted of UPLC-grade water
174 with 0.1% (v/v) FA, and mobile phase B consisted of UPLC-grade ACN with 0.1% (v/v) FA. Peptide elution
175 was performed at 5 μ L/min using the following gradient: i) 2% to 30% mobile phase B in 120 min, ii) ramp
176 to 90% mobile phase B in 1 min. The washing step at 90% mobile phase B lasted 4 min and was followed
177 by an equilibration step at 2% mobile phase B (starting conditions) for 10 min. Ion source parameters
178 were set to 4.5 kV for the ion spray voltage, 25 psi for the curtain gas, 10 psi for nebulizer gas (ion source
179 gas 1), 20 psi for heater gas (ion source gas 2) and 100°C as source temperature.

180 a. Data-Dependent Acquisition

181 For DDA acquisition (a cycle time of 3.3 s), MS1 spectra were collected between 400-1200 m/z for 250
182 ms. The 30 most intense precursor ions with charge states 2-4 that exceeded 300 counts per second
183 were selected for fragmentation, and the corresponding fragmentation MS2 spectra were collected
184 between 100-1500 m/z for 100 ms. After the fragmentation event, the precursor ions were dynamically
185 excluded from reselection for 10 s.

186 b. SWATH 99 Variable windows

187 For SWATH (a cycle time of 4 s), a 99 variable window acquisition scheme was used (see Supplementary
188 Data x) ³⁰. Briefly, SWATH MS2 spectra were collected in high sensitivity mode from 100-1500 m/z, for
189 37.5 ms. Before each SWATH MS cycle an additional MS1 survey scan in high sensitivity mode was
190 recorded for 250 ms.

191 c. Scanning SWATH GP 1Da

192 A Scanning SWATH beta version was installed by a SCIEX engineer on the 6th of October 2020. Scanning
193 SWATH Q1 calibration was performed by directly infusing a tuning solution (ESI Positive Calibration
194 Solution for the SCIEX X500B System - P/N: 5049910) and by acquiring a pre-built calibration batch
195 (SSCalibration.dab). Afterwards, the calibration was verified by i) running a verification calibration batch
196 and inspect the data in PeakView ii) The name of the SSDrift.cal file located in the API
197 Instrument/Preferences folder should be modified which can be checked by looking at the date of last
198 modified.

199 The gas-phase fractionation approach usually acquired in PACIFIC, were also acquired by Scanning
200 SWATH because this uniquely allows to apply DIA annotation algorithms for library building of
201 subsequent full mass range DIA acquisition. Precursor isolation window was set to 1 m/z and a mass
202 range of 100 m/z was covered in 6 s (average accumulation time per precursor: 59.57 ms). An MS1 scan
203 was included and data was acquired in high resolution mode. Raw data were binned in the precursor
204 dimension into 0.2 m/z bins and Q1 calibration was obtained by running rawSSProcessor.exe.

205 d. Scanning SWATH 5Da

206 The Q1 was calibrated again using the same procedure as described in Scanning SWATH GP 1Da. The
207 precursor isolation window was set to 5 m/z and a mass range of 400-900 m/z was covered in 4 s
208 (average accumulation time per precursor: 37.5 ms). A TOF MS scan was included and data was acquired
209 in high sensitivity mode. Raw data were binned in the precursor dimension into 1 m/z bins and Q1
210 calibration was obtained by running the rawSSProcessor.exe.

211 **Note:** rawSSProcessor.exe automatically initializes in the background after a Scanning SWATH run is
212 acquired. We decided not to run the rawSSProcessor.exe on the acquisition desktop itself because we
213 observed LC driver connection issues when rawSSProcessor.exe was running in the background.

214
215 **3) Thermo Orbitrap QE HF-X (Nano flow LC)**

216 A Thermo Orbitrap QE HF-X (Thermo Fisher Scientific, Waltham, Massachusetts, United States) was
217 coupled to an UltiMate 3000 LC-system (NCS-3500RS Nano/Cap System, Thermo Fisher Scientific).
218 Peptides were separated on an Acclaim PepMap C18 column (id 75 µm, length 50 cm, particle size 2
219 µm, Thermo Fisher Scientific ref 164942) at a flow rate of 350 nL/min by means of trap-elute injection
220 (Acclaim PepMap C18, id. 300 µm x 5mm) after 3min desalting on a nano-trap cartridge (id.300 µm,
221 length 5mm, Thermo Fisher Scientific ref 160454).

222 Mobile phase A consisted of UPLC-grade water with 0.1% (v/v) FA, and mobile phase B consisted of
223 UPLC-grade ACN with 0.1% (v/v) FA. Peptide elution was performed at 350 nL/min using the following
224 gradient: i) 2% to 30% mobile phase B in 120 min, ii) ramp to 90% mobile phase B in 1 min. The washing
225 step at 90% mobile phase B lasted 4 min and was followed by an equilibration step at 2% mobile phase
226 B (starting conditions) for 21 min.

227 a. Data-Dependent Acquisition

228 The data-dependent acquisition runs on the Q Exactive HF-X were acquired with MS survey scans (350-
229 1400 m/z) at a resolution of 60,000, and an AGC target of 3e6. The 12 most intense precursor ions,
230 were selected for fragmentation by high-energy collision-induced dissociation, and the resulting
231 fragments were analyzed at a resolution of 15,000 using an AGC target of 1e5 and a maximum fill time
232 of 22 ms. Dynamic exclusion was used within 30 s to prevent repetitive selection of the same peptide.

233 b. Narrow Window Gas-Phase fractionation (GP) DIA

234 Narrow-window GP-DIA data was acquired as described by Searle et al ²⁹. Briefly, 6 GP runs (400-500,
235 ..., 900-1000 m/z) using staggered 4m/z DIA spectra (4m/z precursor isolation windows at 30,000
236 resolution, AGC target 1e6, maximum inject time 60ms, NCE 27, +3H assumed charge state) were
237 acquired using an overlapping window pattern, described by Pino et al ³¹. In each run, full MS scans
238 matching each part of the fractionated mass range (i.e., either 395-505, 495-605, 595-705, 695-805,
239 795-905, or 895-1005 m/z), acquired at a resolution of 60,000 using an AGC target of 1e6 and a
240 maximum inject time of 60ms, were interspersed every 25 MS/MS spectra.

241 c. All Ion Fragmentation (AIF)

242 The AIF DIA data was acquired using a staggered pattern of 75x8 m/z isolation windows over the mass
243 range 400-1000 m/z as described by Pino et al ³¹. DIA MS/MS scans were acquired at 15,000 resolution,
244 with an AGC target of 1e6, a maximum inject time 20ms, and a NCE of 27. Full MS scans over the range
245 390-1010 m/z at 60,000 resolution, AGC target 1e6, maximum inject time 60 ms were interspersed
246 every 75 MS/MS spectra.

247

248 **4) Waters Synapt G2-Si (Capillary flow LC)**

249 An M-class LC system (Waters Corporation, Milford, MA) was equipped with a 1.7 μ m CSH 130 C18 300
250 μ m x 100 mm column, operating at 5 μ L/min with a column temperature of 55 °C. Mobile phase A was
251 UPLC-grade water containing 0.1% (v/v) FA and 3% DMSO, mobile phase B was ACN containing 0.1%
252 (v/v) FA. Peptides were separated using a linear gradient of 3–30% mobile phase B over 120 minutes.
253 All experiments were conducted on a Synapt G2-Si mass spectrometer (Waters Corporation, Wilmslow,
254 UK). The ESI Low Flow probe capillary voltage was 3 kV, sampling cone 60 V, source offset 60 V, source
255 temperature 80 °C, desolvation temperature 350 °C, cone gas 80 L/hr, desolvation gas 350 L/hr, and
256 nebulizer pressure 2.5 bar. A lock mass reference signal of GluFibrinopeptide B (m/z 785.8426) was
257 sampled every 30 s.

258 a. HD-DDA

259 Data was acquired according to Helm et al. with minor adaptations ³². Briefly, in data-dependent mode,
260 the MS automatically switches between MS survey and MS/MS scans based upon a set of switching
261 criteria, including ion intensity and charge state. Full scan MS and MS/MS spectra (m/z 50 - 5000) were
262 acquired in sensitivity mode. MS survey spectra were acquired using a fixed acquisition time of 250 ms
263 and the ions present in each scan were monitored for the following criteria: more than 3000

264 intensity/sec and only 2,3,4,5+ charge states. Once criteria were satisfied, the precursor ion isolation
265 width of the quadrupole was set to 1.0 Th around each precursor sequentially. Tandem mass spectra of
266 up to 12 precursors were generated in the trapping region of the ion mobility cell by using a collisional
267 energy ramp from 6/9 V (low mass 50 Da, start/end) to up to 147/183 V (high mass 5000 Da, start/end),
268 with actual values applied dependent upon the precursor m/z. The MS2 scan time was set to 100 ms
269 and the “TIC stop” parameter was set to 100,000 intensity/s allowing a maximum accumulation time of
270 300 ms (i.e. up to three tandem MS spectra of the same precursor). IMS wave velocity was ramped from
271 2400 m/s to 450 m/s (start to end) and the pusher/ion mobility synchronized for singly charged
272 fragment ions in MS/MS spectra, with up to 85% duty cycle efficiency.

273 b. UDMS^e

274 Data was acquired according to Distler et al. with minor adaptations³³. Briefly, Two data functions were
275 acquired over a mass range of m/z 50 to 2000 in alternating mode, differing only in the collision energy
276 applied to the gas cell. In low-energy MS1 mode, data was collected at a constant gas cell collision
277 energy of 4 eV. In elevated energy MS2 mode, the gas cell collision energy was ramped from 10 to 60
278 eV according to a collision energy look up table in function of drift time. The spectral acquisition time in
279 each mode was 0.6 s with a 0.015 s interscan delay.

280

281 **5) Waters Synapt XS (Capillary flow LC)**

282 An M-class LC system (Waters Corporation, Milford, MA) equipped with a 1.7 μm CSH 130 C18 300 μm
283 x 100 mm column, operating at 7 $\mu\text{L}/\text{min}$ with a column temperature of 55 $^{\circ}\text{C}$ was coupled to a Synapt
284 XS quadrupole ion-TOF mass spectrometer (Waters Corporation, Wilmslow, UK) operating at a mass
285 resolution of 30000, FWHM. The ESI Low Flow probe capillary voltage was 1.8 kV, sampling cone 30 V,
286 source offset 4 V, source temperature 100 $^{\circ}\text{C}$, desolvation temperature 300 $^{\circ}\text{C}$, cone gas disabled,
287 desolvation gas 600 L/hr, and nebulizer pressure 3.5 bar was used. The time-of-flight (TOF) mass
288 analyzer of the mass spectrometer was externally calibrated with a NaCl mixture from m/z 50 to 1990.
289 A lock mass reference signal of GluFibrinopeptide B (m/z 785.8426) was sampled every two minutes.
290 Mobile phase A was water containing 0.1% (v/v) FA, while mobile phase B was ACN containing 0.1%
291 (v/v) FA. The peptides were eluted and separated with a gradient of 5–40% mobile phase B over 120
292 minutes.

293 a. HD-DDA

294 In data-dependent mode, the MS instrument automatically switches between MS survey and MS/MS
295 scans based upon a set of switching criteria, including ion intensity and charge state. Full scan MS and
296 MS/MS spectra (m/z 50 - 5000) were acquired in resolution mode. MS survey spectra were acquired
297 using a fixed acquisition time of 200 ms and the ions present in each scan were monitored for criteria
298 intensity more than 5000 intensity/sec and 2,3,4+ charge states. Once criteria were satisfied, the
299 precursor ion isolation width of the quadrupole was set to 1.0 Th around each precursor sequentially.
300 Tandem mass spectra of up to 15 precursors were generated in the trapping region of the ion mobility
301 cell by using a collisional energy ramp from 6/9 V (low mass 50 Da, start/end) to up to 147/183 V (high
302 mass 5000 Da, start/end), with actual values applied dependent upon the precursor m/z. The MS2 scan
303 time was set to 70 ms and the “TIC stop” parameter was set to 100,000 intensity/s allowing a maximum
304 accumulation time of 100 ms (i.e. up to two tandem MS spectra of the same precursor). IMS wave
305 velocity was ramped from 2450 m/s to 550 m/s (start to end) and the pusher/ion mobility synchronized
306 for singly charged fragment ions in MS/MS spectra, with up to 85% duty cycle efficiency.

307 b. SONAR GP

308 As for Scanning SWATH, the GP fractionation approach which is usually acquired in PACIFIC was also
309 analysed by SONAR purely to extend the peptide library. Therefore, the mass scale from m/z 400 to
310 1200 was divided into 100 Da sections, thus requiring 8 injections for each sample. The quadrupole was
311 continuously scanned from the start mass to end mass of each section and a transmission window of 4
312 Da was used. In low-energy MS1 mode, data were collected at constant gas cell collision energy of 6 eV.
313 In elevated energy MS2 mode, the gas cell collision energy was ramped with values calculated from the
314 start and end mass of the 100 Da mass range being scanned by the quadrupole, and are shown in
315 **Supplementary Table 5**. The spectral acquisition time in each mode was 0.5 s with a 0.02 s interscan
316 delay.

317 c. SONAR

318 The quadrupole was continuously scanned between m/z 400 to 900, with a quadrupole transmission
319 width of ~ 20 Da. Two data functions are acquired in an alternating mode, differing only in the collision
320 energy applied to the gas cell. In low-energy MS1 mode, data were collected at constant gas cell collision
321 energy of 6 eV. In elevated energy MS2 mode, the gas cell collision energy was ramped from 16 to 36
322 eV (per unit charge). The spectral acquisition time in each mode was 0.5 s with a 0.02 s interscan delay.

323

324 **6) Bruker TimsTOF Pro (Nano flow LC)**

325 An Acquity UPLC M-Class System (Waters Corporation) was fitted with a nanoEase™ M/Z Symmetry C18
326 trap column (100Å, 5 μm , 180 μm x 20 mm) and a nanoEase™ M/Z HSS C18 T3 Column (100Å, 1.8 μm ,
327 75 μm x 250 mm, both from Waters Corporation). The sample was loaded onto the trap column in 2min
328 at 5 $\mu\text{l}/\text{min}$ in 94% mobile phase A and 6% mobile phase B. Mobile phase A is UPLC-grade water with
329 0.1% FA, while mobile phase B is 80% ACN with 0.1% FA. The Acquity UPLC M-Class system was coupled
330 online to a TimsTOF Pro via a CaptiveSpray nano-electrospray ion source (Bruker Daltonics, Bremen,
331 Germany), with an ion transfer capillary temperature at 180°C. Liquid chromatography was performed
332 at 40 °C and with a constant flow of 400 nL/min. Peptides were separated using a linear gradient of
333 2–30% mobile phase B over 120 minutes. The TimsTOF Pro elution voltages were calibrated linearly to
334 obtain reduced ion mobility coefficients ($1/K_0$) using three selected ions of the Agilent ESI-L Tuning Mix
335 (m/z , $1/K_0$: 622.0289 Th, 0.9848 Vs cm^{-2} ; 922.0097 Th, 1.1895 Vs cm^{-2} ; 1222.9906 Th, 1.3820 Vs cm^{-2}).

336 a. PASEF

337 Parallel Accumulation–Serial Fragmentation DDA (PASEF) was used to select precursor ions for
338 fragmentation with 1 TIMS-MS scan and 10 PASEF MS/MS scans, as described by Meier et al. in 2018³⁴.
339 The TIMS-MS survey scan was acquired between 0.70 - 1.45 V.s/ cm^2 and 100 - 1700 m/z with a ramp
340 time of 166 ms. The 10 PASEF scans contained on average 12 MS/MS scans per PASEF scan with a
341 collision energy of 10 eV. Precursors with 1 – 5 charges were selected with the target value set to 20
342 000 a.u and intensity threshold to 2500 a.u. Precursors were dynamically excluded for 0.4 min.

343 b. diaPASEF

344 The diaPASEF method was implemented as described by Meier et al. in 2019¹⁹. The DIA parameters
345 that define the windows can be found in supplementary material. The DIA range was set to 400-1200
346 m/z with 16 diaPASEF scans of 25 m/z isolation windows, including an overlap of 1 Da. Each diaPASEF
347 scan consisted of two steps (measuring two 25 Da intervals), with each step spanning an IMS range of
348 0.3 V.s/ cm^2 . The lower IMS value increased linear from 0.6 to 0.834375 for the diaPASEF scans. The
349 TIMS-MS scan was identical to the PASEF method.

350 Data Records

351 Data record 1. The mass spectrometry DDA and DIA-MS proteomics data including instrument raw files
352 (.wiff, .raw, .d) have been deposited to the ProteomeXchange Consortium via the PRIDE partner
353 repository with the dataset identifier, PXD028735^{35,36}. For every instrument, a separate Sample and
354 Data Relationship File (SDRF) and an Investigation Description File (IDF) have been uploaded to
355 ProteomeXchange. Both the SDRF and IDF file formats are relatively new in Proteomics and were
356 developed in a collaboration between EuBIC and the Proteomics Standards Initiative (PSI)^{37,38}. These
357 files are used to annotate the sample metadata and link the metadata to the corresponding data file(s)
358 and thus will improve the reproducibility and reanalysis of this comprehensive benchmark dataset.
359 Reviewer account details are username: reviewer_pxd028735bi.ac.uk and password: NjQ7Nj82.

360 Data record 2. The AutoQC data analysed in Skyline is available from Panorama Public with the link
361 <https://panoramaweb.org/LFQBenchmark.url>.

362

363 Technical Validation

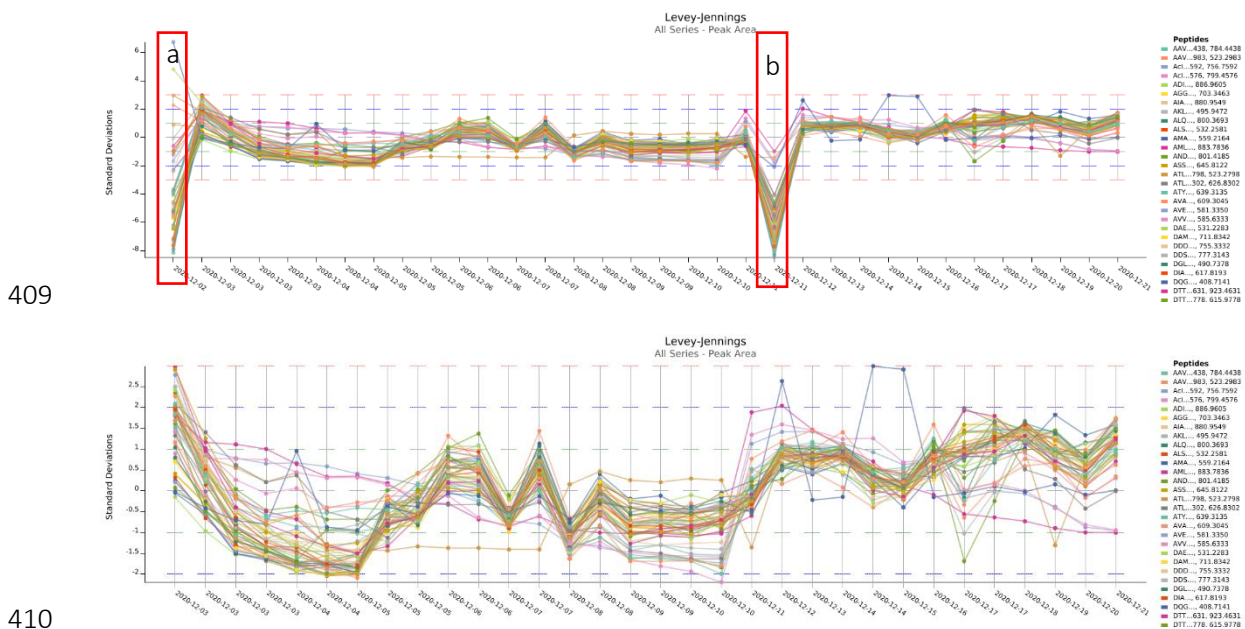
364 We continuously performed system suitability procedures to monitor LC-MS/MS performance in a
365 longitudinal fashion. Therefore, we ran an AutoQC complex lysate, i.e. a commercial E.coli protein digest
366 extract, every 3 to 4 samples over all acquired runs on all instruments. All the AutoQC samples were
367 acquired in DDA on each LC-MS/MS instrument, except for the Synapt XS, Orbitrap QE-HF and the
368 timsTOF Pro, where incidentally DDA and DIA acquisitions were alternated. The same mobile phase A
369 and B composition as for the benchmark samples was used as for the benchmark samples, but the
370 gradient applied was modified to reduce the time required to acquire the complete sample batch: linear
371 3-40% B in 60 minutes, up to 85% B in 2 minutes, isocratic at 85% B for 7 minutes, down to 3% B in 1
372 minute and isocratic at 3% B for 10 minutes. Note that the timsTOF Pro AutoQC samples were acquired
373 using the 120min gradient similar to the actual hybrid proteome samples.

374 System suitability assessment was performed by monitoring peptide-identification free metrics (i.e.
375 retention time, peak area, mass accuracy, etc.) extracted with the vendor neutral Panorama AutoQC
376 framework^{39,40}. To isolate a set of peptides that can be used for this, triplicate AutoQC samples acquired
377 on each instrument were peak picked using MSConvert (version 3.0.20070) and the corresponding .MGF
378 files were searched against an E.coli FASTA database using Mascot Daemon (v2.7). The searches were
379 performed with following settings: (i) 20 ppm peptide mass tolerance, (ii) 50 ppm fragment mass
380 tolerance and (iii) two allowed missed cleavages. The peptide and protein identification results were
381 exported as Mascot .DAT file and imported into Skyline Daily (version 21.1.1.160). The five highest
382 ranked proteins were retained in the target list and after importing one of the AutoQC .raw files, we
383 manually verified and removed each precursor with co-eluting peptides and low MS1 signal intensity
384 before a Skyline file was saved as template file. Finally, a configuration file for each setup was created
385 with the AutoQC Loader software (version 21.1.0.158) which leads to the automatic import of every
386 sample, with the pattern "AutoQC" in the .raw file or folder structure, in the Skyline template. The data
387 and skyline reports were published to the PanoramaWeb folder "U of Ghent Pharma Biotech Lab -
388 LFQBenchmark across Instrument Platforms" containing six subfolders for each instrumental platform.

389 For each instrument, peak area, retention time and mass accuracy were manually checked by plotting
390 these metrics in Levey-Jennings plots. For almost every instrument a few outliers were detected, as can
391 be expected on a dataset of over 600 LCMS runs. Fortunately, most of these can be explained by
392 inspecting the raw data and by personal communication with the technicians acquiring the respective
393 datasets. **Figure 2** illustrates one such case. More specifically, two AutoQC samples display a near-

394 complete loss in peak area in the TripleTOF6600+ DDA dataset. Indeed, these were caused by (a) a
395 wrong vial in the sample tray and (b) an empty vial. When these two samples are removed from the QC
396 plot, a more coherent perspective on the variation in standard deviation is seen in the Levey-Jennings
397 plot. Other instances that we have already found include (i) a significant shift in standard deviation in
398 peak area reported for both the Orbitrap, timsTOF Pro and Synapt XS dataset because AutoQC samples
399 were incidentally acquired in two different acquisition methodologies, i.e. in DDA and DIA; (ii) For the
400 timsTOF Pro, a drift in retention time was seen, indicating LC related technical variation which could
401 have been caused by e.g. too short column equilibration times; (iii) In the Orbitrap QE HF-X AutoQC data
402 one peptide (EEAIK) was undetectable in all the AutoQC samples acquired in AIF. Manual inspection of
403 the acquisition in Skyline (easily accessible through the Panorama QC pipeline) surfaced that it fell out
404 of the precursor m/z range (351.7053) that was acquired.

405 As expected in such a massive MS proteomics experiments, it seems that for every instrument some
406 outliers were recorded, most of which have explanations common to the field. Above all, this
407 demonstrates the necessity of a performant system suitability workflow to increase the reproducibility
408 and quality of LC-MS/MS proteomics datasets⁴¹.



410
411 **Figure 2. Levey-Jennings plot of the standard deviation in peak area for 50 selected precursors acquired**
412 **in DDA with the TripleTOF6600+. The upper chart shows two distinct outliers, acquired respectively on**
413 **the 2nd and 12th of December (red boxes). Manual inspection of the data shows these were caused by**
414 **(a) a wrong vial in the sample tray and (b) an empty vial. When these two samples are excluded from**
415 **the Levey-Jennings plot (lower chart), a significant drop in standard deviation over the time period of**
416 **data acquisition is seen.**

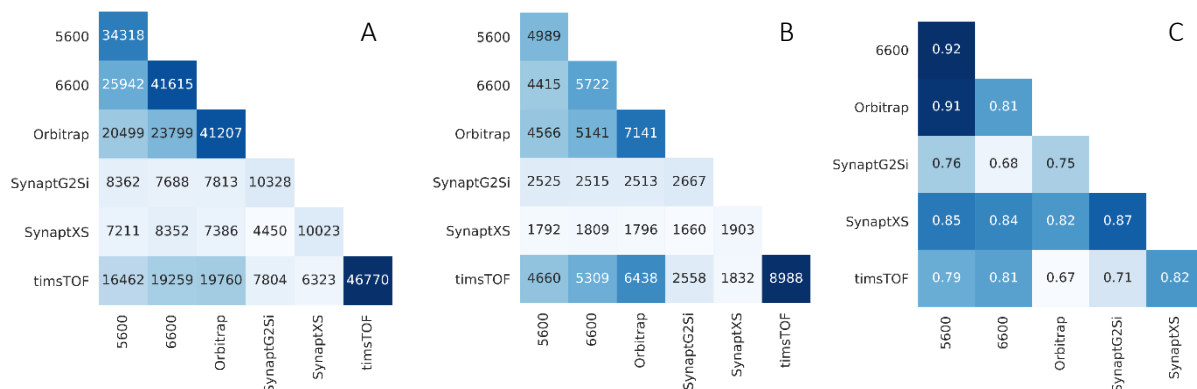
417 Usage Notes

418 A comprehensive dataset such as the one presented here is inspired by a bioinformatics need to cope
419 with the recent expansion of novel acquisition strategies and data dimensions. Apart from being a
420 repository for validating both performance and compatibility of (new) bioinformatic pipelines, it can
421 serve as a reference for general proteomics courses (e.g. Skyline tutorials, SWATH/DIA course) and be
422 applied for training and validating machine/deep learning algorithms. As such, it is intended to facilitate
423 our understanding of the impact of instrumentation on the perspective that is generated on protein
424 biology and to a larger extent to help unify the field of proteomics.

425 To demonstrate the applicability of this data repository, we assessed the impact of instrumentation on
 426 the most conventional data format acquired by all instruments i.e. DDA. Therefore, for every
 427 instrumental platform, triplicate Human, Yeast and E.coli DDA runs were peak picked with MSConvert
 428 (version 3.0.21285) and exported as Mascot .MGF file. This software was chosen as it is vendor-
 429 independent and contains each vendor's implementation for peak picking, with the exception of UNIFI
 430 i.e. Waters. Therefore, Progenesis QI for Proteomics (version 4.2.7207) was used for the Waters DDA
 431 data. By doing so, the MS1 precursor space was aligned in the retention time dimension before peak
 432 picking. Subsequently, all MS/MS spectra were exported as .MGF file for peptide identification. A
 433 standard search with carbamidomethylation of cysteine as fixed modification was performed using a
 434 database containing the Human, Yeast and E.coli protein sequences (downloaded from Uniprot on
 435 19/01/2021) and with following parameters: i) mass error tolerances for the precursor ions and the
 436 fragment ions were set at 20 ppm and 50 ppm, respectively; ii) enzyme specificity was set to trypsin,
 437 allowing up to one missed cleavage. Next, the results were exported as .DAT file and imported into
 438 Skyline to create a non-redundant spectral library with BiblioSpec. Afterwards, the .BLIB file was
 439 converted to a .dlib and .msp file format respectively using EncyclopeDIA²⁹. The resulting .msp file was
 440 converted using a Python conversion tool (speclib_to_mgf.py) built-in MS²PIP, to create a peptide
 441 record file (.PEPREC) and .MGF file. Next the proportion (amount) of peptide identifications overlapping
 442 between the different instruments was assessed using a custom Python script (**Figure 3A**).

443 Since each instrument analysed the same commercial protein digests, a large overlap in peptide
 444 identifications would be expected. However, while the timsTOF Pro, TripleTOF 6600+ and Orbitrap QE
 445 HF-X roughly identify a similar number of peptide sequences (approximately 40,000), the overlap in
 446 identified sequences is in the order of 50%, with that overlap between the Orbitrap QE HF-X and the
 447 TripleTOF5600 and 6600+ being overall 10% higher than the overlap with the timsTOF Pro. Fortunately,
 448 at the protein level these differences tend to flatten (**Figure 3B**). Still, such remarkable findings require
 449 follow-up analyses that can resolve the impact of the differences in instrumental design on the kind of
 450 peptides that can be detected. In the process providing a deeper understanding of the impact on the
 451 underlying biology that can be studied on different instrument designs can be provided.

452 In pursuit of one such follow-up analysis, we assessed the differences in the fragmentation process
 453 between the instruments (all of which are beam-type CID) by first mutually mapping the MS2 intensities
 454 and plotting their median Pearson correlation coefficients (PCC) (**Figure 3C**). By definition, this only plots
 455 commonly found peptides. Indeed, the largest differences in fragment intensities are found between
 456 timsTOF Pro and Orbitrap QE HF-X and this could therefore underlie a difference in performance of the
 457 annotation algorithm used, i.e. Mascot. However, digging into the algorithmic intricacies of Mascot that
 458 could couple fragment intensities to differential identification is outside the scope of this manuscript.

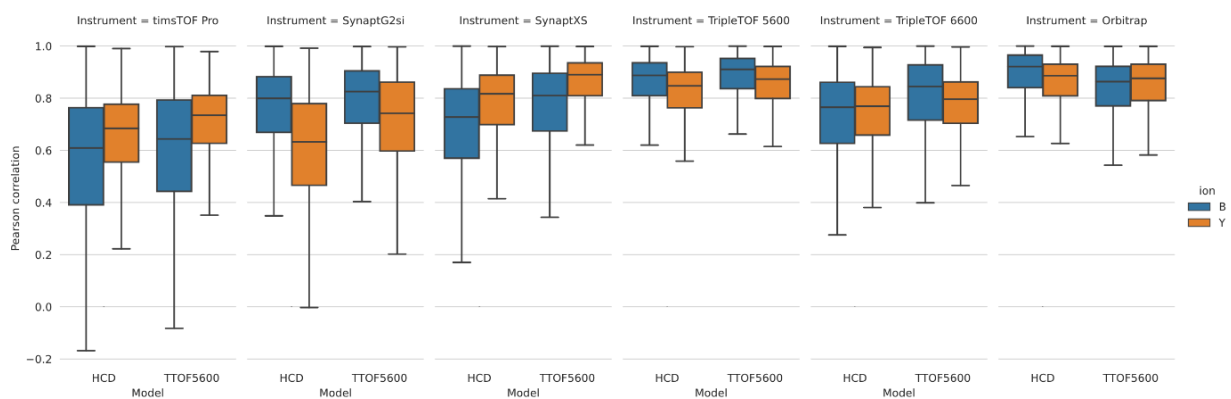


459 **Figure 3. Correlation matrices comparing the DDA data of six different instruments** in number of identified
 460 peptide and protein sequences, and fragmentation similarity expressed as Pearson Correlation

461 Coefficient (PCC). A) Describes the overlap in uniquely identified peptide sequences, while B) shows the
462 overlap in protein sequences from triplicate Human, Yeast and E.coli DDA samples between the six
463 instruments. C) PCC were calculated between the shared identified peptides from the DDA replicates
464 between each instrument. The numbers in each box correspond to the median spectrum PCC between
465 the instrument on the x-axis and the instrument on the y-axis. A dark blue color indicates a higher degree
466 of overlap or higher median PCC.

467 Still, fragment intensities have become increasingly important in proteomics with the introduction of
468 machine learning algorithms that can predict fragment intensities based on a simple peptide sequence
469 input^{24,25,42-44}. Therefore, to attain an even deeper understanding of these fragmentation differences,
470 we calculated the PCC of both b- and y-ions compared to two prediction models (i.e. HCD and TTOF
471 5600) from MS²PIP⁴⁵. This confirms that the orbitrap and triple TOF designs have a very similar
472 fragmentation pattern.

473 Excitingly, this analysis surfaces a deeper insight: instruments containing an ion mobility separation
474 device, i.e. the timsTOF Pro (TIMS) and Synapt series instruments (TWIMS), produced deviating
475 fragment intensities **Figure 4**. Importantly, Waters and Bruker apply IMS in a very different way, forcing
476 a deeper understanding of instrumental architecture and how ion mobility is applied by both vendors
477 in DDA. Briefly, the order of ion manipulation is Q-CID-IMS in the so-called High Definition DDA (HD-
478 DDA) in the Waters series and it is IMS-Q-CID in the Bruker instruments. Therefore, Waters separates
479 the fragment ions in IMS and leverages the efficient charge state separation to synchronize the pusher
480 of the TOF tube with singly charged fragment ions in order to only and nearly one hundred percent
481 efficiently sample the single charged fragments⁴⁶. In the timsTOF Pro on the other hand, the IMS is in
482 fact resolving the precursor ion space before fragmentation, leading to a different selection of the
483 peptide precursor space compared to other devices that do not use IMS (in this way). This most
484 prominently underlies the surprisingly low overlap with e.g. the Orbitrap HFX, yet does not at first sight
485 explain the difference in fragment intensities. However, one possible reason could be the energy
486 dependency of beam-type CID i.e. normalised collision energy, which a few prediction algorithms (Prosit
487 and pDeep3) have included as input feature^{25,47,48}.



488

489 **Figure 4. Boxplots of the Pearson correlation coefficients (PCC) between the MS²PIP predicted (HCD and**
490 **TTOF5600 model) and experimental spectra across the six different LC-MS instruments.**

491 In conclusion, it is clear from this preliminary data usage case on the simple and most commonly applied
492 DDA strategies that a lot of insights on data structure and bias will be generated, especially in light of
493 the more recent DIA strategies. Importantly, each step in the data processing can greatly impact the
494 final outcome and we especially anticipate a renewed interest in (multidimensional) peak picking
495 algorithms. Especially for the latest generation of DIA acquisition methods, large differences in
496 annotation performance are known to exist between annotation tools, despite the optional

497 compatibility that some of these offer. We would therefore like to invite the developers of all current
498 bioinformatics tools, including the very recent peptide intensity predictors e.g. Prosit, pDeep, and
499 DeepDIA to make use of this comprehensive dataset to benchmark their performance for each
500 instrument individually and adapt their algorithms to increase the performance on all. For us it is clear
501 that especially the complementarity of the detectable and annotatable ion space will move the field
502 forward and help researchers make informed decisions on the best acquisition strategies for their
503 application or biological question under investigation.

504 **Code Availability**

505 MS²PIP is open source, licensed under the Apache-2.0 License, and hosted on
506 https://github.com/compomics/ms2pip_c. The Jupyter notebooks used to generate **Figure 3** and **4** are
507 available through Zenodo, under DOI: 10.5281/zenodo.5714380.

508 **Acknowledgements**

509 This research was funded by grants from the Research Foundation Flanders (FWO) awarded to BVP
510 (grant number: 11B4518N), RG (1S50918N), and MD (12E9716N). Hans Vissers is acknowledged for his
511 assistance with data conversion and formatting. This work was supported in part by the French Ministry
512 of Research with the Investissement d’Avenir Infrastructures Nationales en Biologie et Santé program
513 (ProFi, Proteomics French Infrastructure project; ANR-10-INBS-08).

514 **Author contributions**

515 BVP, SW, SD and MD conceived the study, BVP performed the TripleTOF 5600 and 6600+ data
516 acquisition, SD performed the Synapt G2-Si data acquisition, AGP, DB, EM and KC performed the
517 Orbitrap data acquisition, KB performed the timsTOF Pro data acquisition. CH and LG performed the
518 Synapt XS data acquisition. BVP and SD prepared the samples and RG wrote the scripts to generate
519 figures 3 and 4. YPR organised the ProteomeXchange submission. BVP and MD wrote the draft
520 manuscript with contributions from all authors. MD supervised and DD and LM co-supervised the
521 experiment.

522 **Competing interests**

523 Chris Hughes and Lee Gethings are employed by Waters Corporation.

524 **References**

- 525 1. Willems, S. *et al.* Ion-networks: A sparse data format capturing full data integrity of data
526 independent acquisition mass spectrometry. *bioRxiv* (2019) doi:10.1101/726273.
- 527 2. Vaudel, M. *et al.* PeptideShaker enables reanalysis of MS-derived proteomics data sets. *Nat*
528 *Biotechnol* **33**, 22–24 (2015).
- 529 3. Verheggen, K., Martens, L., Berven, F. S., Barsnes, H. & Vaudel, M. Database Search Engines:
530 Paradigms, Challenges and Solutions. in *Advances in Experimental Medicine and Biology* 147–
531 156 (2016). doi:10.1007/978-3-319-41448-5_6.
- 532 4. Geromanos, S. J., Hughes, C., Ciavarini, S., Vissers, J. P. C. & Langridge, J. I. Using ion purity scores
533 for enhancing quantitative accuracy and precision in complex proteomics samples. *Anal. Bioanal.*
534 *Chem.* **404**, 1127–1139 (2012).
- 535 5. Richardson, K. *et al.* A probabilistic framework for peptide and protein quantification from data-
536 dependent and data-independent LC-MS proteomics experiments. *OMICS* **16**, 468–482 (2012).

- 537 6. Li, G. Z. *et al.* Database searching and accounting of multiplexed precursor and product ion
538 spectra from the data independent analysis of simple and complex peptide mixtures. *Proteomics*
539 **9**, 1696–1719 (2009).
- 540 7. Helm, D. *et al.* Ion Mobility Tandem Mass Spectrometry Enhances Performance of Bottom-up
541 Proteomics. *Mol. Cell. Proteomics* **13**, 3709–3715 (2014).
- 542 8. Shliaha, P. V, Bond, N. J., Gatto, L. & Lilley, K. S. Effects of traveling wave ion mobility separation
543 on data independent acquisition in proteomics studies. *J Proteome Res* **12**, 2323–2339 (2013).
- 544 9. Gillet, L. C. *et al.* Targeted data extraction of the MS/MS spectra generated by data-independent
545 acquisition: a new concept for consistent and accurate proteome analysis. *Mol Cell Proteomics*
546 **11**, O111 016717 (2012).
- 547 10. Ludwig, C. *et al.* Data-independent acquisition-based SWATH-MS for quantitative proteomics: a
548 tutorial. *Mol. Syst. Biol.* (2018) doi:10.15252/msb.20178126.
- 549 11. Ting, Y. S. *et al.* Peptide-Centric Proteome Analysis: An Alternative Strategy for the Analysis of
550 Tandem Mass Spectrometry Data. *Mol. Cell. Proteomics* **14**, 2301–7 (2015).
- 551 12. Li, Y. *et al.* Group-DIA: analyzing multiple data-independent acquisition mass spectrometry data
552 files. *Nat. Methods* **12**, 1105–1106 (2015).
- 553 13. Kuharev, J., Navarro, P., Distler, U., Jahn, O. & Tenzer, S. In-depth evaluation of software tools
554 for data-independent acquisition based label-free quantification. *Proteomics* **15**, 3140–3151
555 (2015).
- 556 14. Teleman, J. *et al.* DIANA-algorithmic improvements for analysis of data-independent acquisition
557 MS data. *Bioinformatics* (2015) doi:10.1093/bioinformatics/btu686.
- 558 15. Peckner, R. *et al.* Specter: linear deconvolution for targeted analysis of data-independent
559 acquisition mass spectrometry proteomics. *Nat. Methods* **15**, 371–378 (2018).
- 560 16. Wang, J. *et al.* MSPLIT-DIA: sensitive peptide identification for data-independent acquisition.
561 *Nat. Methods* **12**, 1106–1108 (2015).
- 562 17. Moseley, M. A. *et al.* Scanning Quadrupole Data-Independent Acquisition, Part A: Qualitative and
563 Quantitative Characterization. *J. Proteome Res.* **17**, 770–779 (2018).
- 564 18. Messner, C. B. *et al.* Ultra-fast proteomics with Scanning SWATH. *Nat. Biotechnol.* **39**, 846–854
565 (2021).
- 566 19. Meier, F. *et al.* diaPASEF: parallel accumulation–serial fragmentation combined with data-
567 independent acquisition. *Nat. Methods* **17**, 1229–1236 (2020).
- 568 20. Van Puyvelde, B. *et al.* Removing the Hidden Data Dependency of DIA with Predicted Spectral
569 Libraries. *Proteomics* **20**, e1900306 (2020).
- 570 21. Demichev, V., Messner, C. B., Vernardis, S. I., Lilley, K. S. & Ralser, M. DIA-NN: neural networks
571 and interference correction enable deep proteome coverage in high throughput. *Nat. Methods*
572 **17**, 41–44 (2019).
- 573 22. Bouwmeester, R., Gabriels, R., Van Den Bossche, T., Martens, L. & Degroeve, S. The Age of Data-
574 Driven Proteomics: How Machine Learning Enables Novel Workflows. *Proteomics* (2020)
575 doi:10.1002/pmic.201900351.
- 576 23. C. Silva, A. S., Bouwmeester, R., Martens, L. & Degroeve, S. Accurate peptide fragmentation
577 predictions allow data driven approaches to replace and improve upon proteomics search engine
578 scoring functions. *Bioinformatics* **35**, 5243–5248 (2019).

- 579 24. Zhou, X. X. *et al.* PDeep: Predicting MS/MS Spectra of Peptides with Deep Learning. *Anal. Chem.*
580 (2017) doi:10.1021/acs.analchem.7b02566.
- 581 25. Gessulat, S. *et al.* ProSIT: proteome-wide prediction of peptide tandem mass spectra by deep
582 learning. *Nat. Methods* (2019) doi:10.1038/s41592-019-0426-7.
- 583 26. Bouwmeester, R., Gabriels, R., Hulstaert, N., Martens, L. & Degroeve, S. DeepLC can predict
584 retention times for peptides that carry as-yet unseen modifications. *bioRxiv* 2020.03.28.013003
585 (2021) doi:10.1101/2020.03.28.013003.
- 586 27. Mann, M., Kumar, C., Zeng, W.-F. & Strauss, M. T. Artificial intelligence for proteomics and
587 biomarker discovery. *Cell Syst.* **12**, 759–770 (2021).
- 588 28. Navarro, P. *et al.* A multicenter study benchmarks software tools for label-free proteome
589 quantification. *Nat. Biotechnol.* **34**, 1130–1136 (2016).
- 590 29. Searle, B. C. *et al.* Chromatogram libraries improve peptide detection and quantification by data
591 independent acquisition mass spectrometry. *Nat. Commun.* **9**, 5128 (2018).
- 592 30. Improved Data Quality Using Variable Q1 Window Widths in SWATH[®] Acquisition Data
593 Independent Acquisition on TripleTOF[®] and X-Series QTOF Systems. (2019).
- 594 31. Pino, L. K., Just, S. C., MacCoss, M. J. & Searle, B. C. Acquiring and Analyzing Data Independent
595 Acquisition Proteomics Experiments without Spectrum Libraries. *Mol. Cell. Proteomics* **19**, 1088–
596 1103 (2020).
- 597 32. Helm, D. *et al.* Ion mobility tandem mass spectrometry enhances performance of bottom-up
598 proteomics. *Mol Cell Proteomics* **13**, 3709–3715 (2014).
- 599 33. Distler, U. *et al.* Drift time-specific collision energies enable deep-coverage data-independent
600 acquisition proteomics. *Nat. Methods* **11**, 167–170 (2014).
- 601 34. F, M. *et al.* Online Parallel Accumulation-Serial Fragmentation (PASEF) with a Novel Trapped Ion
602 Mobility Mass Spectrometer. *Mol. Cell. Proteomics* **17**, 2534–2545 (2018).
- 603 35. JA, V. *et al.* 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res.* **44**, D447–
604 D456 (2016).
- 605 36. Perez-Riverol, Y. *et al.* The PRIDE database and related tools and resources in 2019: Improving
606 support for quantification data. *Nucleic Acids Res.* (2019) doi:10.1093/nar/gky1106.
- 607 37. Dai, C. *et al.* A proteomics sample metadata representation for multiomics integration and big
608 data analysis. *Nat. Commun.* 2021 121 **12**, 1–8 (2021).
- 609 38. Bittremieux, W. *et al.* The European Bioinformatics Community for Mass Spectrometry (EuBIC-
610 MS): an open community for bioinformatics training and research. *Rapid Commun. Mass*
611 *Spectrom.* e9087 (2021) doi:10.1002/RCM.9087.
- 612 39. Bereman, M. S. *et al.* An Automated Pipeline to Monitor System Performance in Liquid
613 Chromatography-Tandem Mass Spectrometry Proteomic Experiments. *J Proteome Res* **15**, 4763–
614 4769 (2016).
- 615 40. Sharma, V. *et al.* Panorama: A Targeted Proteomics Knowledge Base. *J. Proteome Res.* **13**, 4205–
616 4210 (2014).
- 617 41. Bereman, M. S. Tools for monitoring system suitability in LC MS/MS centric proteomic
618 experiments. *Proteomics* **15**, 891–902 (2015).
- 619 42. Degroeve, S. & Martens, L. MS2PIP: a tool for MS/MS peak intensity prediction. *Bioinformatics*

620 **29**, 3199–3203 (2013).

621 43. Tiwary, S. *et al.* High-quality MS/MS spectrum prediction for data-dependent and data-
622 independent acquisition data analysis. *Nat. Methods* 2019 166 **16**, 519–525 (2019).

623 44. Yang, Y. *et al.* In silico spectral libraries by deep learning facilitate data-independent acquisition
624 proteomics. *Nat. Commun.* 2020 111 **11**, 1–11 (2020).

625 45. Gabriels, R., Martens, L. & Degroeve, S. Updated MS²PIP web server delivers fast and accurate
626 MS² peak intensity prediction for multiple fragmentation methods, instruments and labeling
627 techniques. *Nucleic Acids Res.* **47**, W295–W299 (2019).

628 46. Helm, D. *et al.* Ion Mobility Tandem Mass Spectrometry Enhances Performance of Bottom-up
629 Proteomics. *Mol. Cell. Proteomics* **13**, 3709–3715 (2014).

630 47. JK, D., AF, P. & JR, Y. Energy dependence of HCD on peptide fragmentation: stepped collisional
631 energy finds the sweet spot. *J. Am. Soc. Mass Spectrom.* **24**, 1690–1699 (2013).

632 48. Tarn, C. & Zeng, W.-F. pDeep3: Toward More Accurate Spectrum Prediction with Fast Few-Shot
633 Learning. *Anal. Chem.* **93**, 5815–5822 (2021).

634

635

636

637

638

639

640

641

642

643